Reinforcement Learning



final level limited by teacher

"Reinforcement learning"



student/scientist (tries out things)

final level: unlimited (?)

Reinforcement learning



fully observed vs. partially observed "state" of the environment

Self-driving cars, robotics:

Observe immediate environment & move

Games:

Observe board & place stone

Observe video screen & move player

Challenge: the "correct" action is not known! Therefore: no supervised learning!

Reward will be rare (or decided only at end)

Use **reinforcement learning**:

Training a network to produce actions based on rare rewards (instead of being told the 'correct' action!)

Challenge: We could use the final reward to define a cost function, but we cannot know how the environment reacts to a proposed change of the actions that were taken!

(unless we have a model of the environment)

Reinforcement Learning: Basic setting





RL-environment





=REINFORCE (Williams 1992): The simplest **model-free** general reinforcement learning technique



Basic idea: Use probabilistic action choice. If the reward at the end turns out to be high, make **all** the actions in this sequence **more likely** (otherwise do the opposite)

This will also sometimes reinforce 'bad' actions, but since they occur more likely in trajectories with low reward, the net effect will still be to suppress them!



RL-environment

Policy: $\pi_{\theta}(a_t|s_t)$ – probability to pick action a_t given observed state s_t at time t



RL-environment



RL-environment

Probabilistic policy:

Probability to take action a, given the current state s

 $\pi_{\theta}(a|s)$

parameters of the network



Environment: makes (possibly stochastic) transition to a new state s', and possibly gives a reward r

Transition function P(s'|s, a)

Probability for having a certain trajectory of actions and states: product over time steps

$$P_{\theta}(\tau) = \Pi_t P(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$$

trajectory: $au = (\mathbf{a}, \mathbf{s})$

 $\mathbf{a} = \dot{a}_0, \dot{a}_1, a_2, \dots$ $\mathbf{s} = s_1, s_2, \dots$ (state 0 is fixed)

Expected overall reward (='return'): sum over all trajectories

$$\bar{R} = E[R] = \sum_{\tau} P_{\theta}(\tau) R(\tau) - \operatorname{return \ for \ this \ sequence \ (sum \ over \ individual \ rewards \ r \ for \ all \ times)}$$

sum over all actions at all times and over all states at all times >0

$$\sum_{\tau} \ldots = \sum_{a_0, a_1, a_2, \ldots, s_1, s_2, \ldots} \ldots$$

Try to maximize expected return by changing parameters of policy: $2\overline{D}$

$$\frac{\partial R}{\partial \theta} = ?$$



Main formula of policy gradient method:

$$\frac{\partial \bar{R}}{\partial \theta} = \sum_{t} E[R \frac{\partial \ln \pi_{\theta}(a_t | s_t)}{\partial \theta}]$$

Stochastic gradient descent:

 $\Delta \theta = \eta \frac{\partial R}{\partial \theta} \quad \text{where } E[\dots] \text{ is approximated via the value for one trajectory (or a batch)}$

$$\frac{\partial \bar{R}}{\partial \theta} = \sum_{t} E[R \frac{\partial \ln \pi_{\theta}(a_t | s_t)}{\partial \theta}]$$

Increase the probability of all action choices in the given sequence, depending on size of return R. Even if R>0 always, due to normalization of probabilities this will tend to suppress the action choices in sequences with lower-than-average returns.

Abbreviation: $G_{k} = \frac{\partial \ln P_{\theta}(\tau)}{\partial \theta_{k}} = \sum_{t} \frac{\partial \ln \pi_{\theta}(a_{t}|s_{t})}{\partial \theta_{k}}$ $\frac{\partial \bar{R}}{\partial \theta_{k}} = E[RG_{k}]$

Policy Gradient: reward baseline

Challenge: fluctuations of estimate for return gradient can be huge. Things improve if one subtracts a constant baseline from the return.

$$\frac{\partial \bar{R}}{\partial \theta} = \sum_{t} E[(R - b) \frac{\partial \ln \pi_{\theta}(a_t | s_t)}{\partial \theta}]$$
$$= E[(R - b)G]$$

This is the same as before. Proof:

$$E[G_k] = \sum_{\tau} P_{\theta}(\tau) \frac{\partial \ln P_{\theta}(\tau)}{\partial \theta_k} = \frac{\partial}{\partial \theta_k} \sum_{\tau} P_{\theta}(\tau) = 0$$

However, the variance of the fluctuating random variable (R-b)G is different, and can be smaller (depending on the value of b)!

Optimal baseline

$$X_k = (R - b_k)G_k$$
$$\operatorname{Var}[X_k] = E[X_k^2] - E[X_k]^2 = \min$$
$$\frac{\partial \operatorname{Var}[X_k]}{\partial b_k} = 0$$

$$b_k = \frac{E[G_k^2 R]}{E[G_k^2]}$$

$$G_k = \frac{\partial \ln P_\theta(\tau)}{\partial \theta_k}$$
$$\Delta \theta_k = -\eta E[G_k(R - b_k)]$$

For more in-depth treatment, see David Silver's course on reinforcement learning (University College London):

http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html

A random walk, where the probability to go "up" is determined by the policy, and where the return is given by the final position (ideal strategy: always go up!)

(Note: this policy does not even depend on the current state)



A random walk, where the probability to go "up" is determined by the policy, and where the return is given by the final position (ideal strategy: always go up!)

(Note: this policy does not even depend on the current state)

policy $\pi_{\theta}(up) = \frac{1}{1 + e^{-\theta}}$ return R = x(T)**RL update** $\Delta \theta = \eta \sum_{t} \left\langle R \frac{\partial \ln \pi_{\theta}(a_t)}{\partial \theta} \right\rangle$ $a_t = up \text{ or down}$ $\frac{\partial \ln \pi_{\theta}(a_t)}{\partial \theta} = \pm e^{-\theta} \pi_{\theta}(a_t) = \pm (1 - \pi_{\theta}(a_t)) = \frac{1 - \pi_{\theta}(\text{up}) \text{ for up}}{-\pi_{\theta}(\text{up}) \text{ for down}}$ $\sum_{t} \frac{\partial \ln \pi_{\theta}(a_{t})}{\partial \theta} = N_{up} - N \pi_{\theta}(up)$ N=number of time steps

return
$$R = x(T) = N_{up} - N_{down} = 2N_{up} - N$$

RL update $\Delta \theta = \eta \sum_{t} \left\langle R \frac{\partial \ln \pi_{\theta}(a_{t})}{\partial \theta} \right\rangle$
 $a_{t} = up \text{ or down}$
 $\left\langle R \sum_{t} \frac{\partial \ln \pi_{\theta}(a_{t})}{\partial \theta} \right\rangle = 2 \left\langle (N_{up} - \frac{N}{2})(N_{up} - \bar{N}_{up}) \right\rangle$
(general analytical expression for average update, rare)
Initially, when $\pi_{\theta}(up) = \frac{1}{2}$:
 $\Delta \theta = 2\eta \left\langle (N_{up} - \frac{N}{2})^{2} \right\rangle = 2\eta \text{Var}(N_{up}) = \eta \frac{N}{2} > 0$
(binomial distribution!)

In general:

$$\left\langle R \sum_{t} \frac{\partial \ln \pi_{\theta}(a_{t})}{\partial \theta} \right\rangle = 2 \left\langle \left(N_{\rm up} - \frac{N}{2} \right) (N_{\rm up} - \bar{N}_{\rm up}) \right\rangle$$
$$= 2 \left\langle \left(\left(N_{\rm up} - \bar{N}_{\rm up} \right) + \left(\bar{N}_{\rm up} - \frac{N}{2} \right) \right) (N_{\rm up} - \bar{N}_{\rm up}) \right\rangle$$
$$= 2 \mathrm{Var} N_{\rm up} + 2 (\bar{N}_{\rm up} - \frac{N}{2}) \left\langle N_{\rm up} - \bar{N}_{\rm up} \right\rangle$$
$$= 2 \mathrm{Var} N_{\rm up} = 2 N \pi_{\theta} (\mathrm{up}) (1 - \pi_{\theta} (\mathrm{up})) \quad \text{(general analytication for the second of the$$



(general analytical expression for average update, fully simplified, extremely rare)



(This plot for N=100 time steps in a trajectory; eta=0.001)

Spread of the update step



Optimal baseline suppresses spread!



Note: Many update steps reduce relative spread

M = number of update steps





$$\frac{\sqrt{\mathrm{Var}\Delta X}}{\langle \Delta X \rangle} \sim \frac{1}{\sqrt{M}}$$

Implement the RL update including the optimal baseline and run some stochastic learning attempts. Can you observe the improvement over the no-baseline results shown here?

Note: You do not need to simulate the individual random walk trajectories, just exploit the binomial distribution.

The second-simplest RL example



See code on website: "SimpleRL_WalkerTarget"

output = action probabilities (softmax) policy $\pi_{\theta}(a|s)$



Policy gradient: all the steps

Obtain one "trajectory":



Policy gradient: all the steps

For each trajectory:



RL in keras: categorical cross-entropy trick

output = action probabilities (softmax) $\pi_{\theta}(a|s)$



categorical cross-entropy

$$C = -\sum_{\substack{a \\ a \\ \text{desired} \\ \text{distribution}}} \frac{\text{distr. from net}}{\pi_{\theta}(a|s)}$$

Set

$$P(a) = R$$

for a=action that was taken

 $P(a) = 0 \label{eq:particular}$ for all other actions a

$$\Delta \theta = -\eta \frac{\partial C}{\partial \theta}$$

implements policy gradient

RL in keras: categorical cross-entropy trick

Encountered N states (during repeated runs)

After setting categorical cross-entropy as cost function, just use the following simple line to implement policy gradient: array N x state-size

Here **desired_outputs[j,a]=R** for the state numbered **j**, if action **a** was taken during a run that gave overall return **R**



Among the major board games, "Go" was not yet played on a superhuman level by any program (very large state space on a 19x19 board!)

alpha-Go beat the world's best player in 2017

First: try to learn from human expert players

sampled state-action pairs (*s*, *a*), using stochastic gradient ascent to maximize the likelihood of the human move *a* selected in state *s*

$$\Delta \sigma \propto \frac{\partial \log p_{\sigma}(a \mid s)}{\partial \sigma}$$

We trained a 13-layer policy network, which we call the SL policy network, from 30 million positions from the KGS Go Server. The net-

Silver et al., "Mastering the game of Go with deep neural networks and tree search" (Google Deepmind team), Nature, January 2016

Second: use policy gradient RL on games played against previous versions of the program

to the current policy. We use a reward function r(s) that is zero for all non-terminal time steps t < T. The outcome $z_t = \pm r(s_T)$ is the terminal reward at the end of the game from the perspective of the current player at time step t: +1 for winning and -1 for losing. Weights are then updated at each time step t by stochastic gradient ascent in the direction that maximizes expected outcome²⁵

$$\Delta
ho \propto rac{\partial \log p_{
ho}(a_t | s_t)}{\partial
ho} z_t$$

Silver et al., "Mastering the game of Go with deep neural networks and tree search" (Google Deepmind team), Nature, January 2016



*Note: beyond policygradient type methods, this also includes another algorithm, called Monte Carlo Tree Search

Silver et al., "Mastering the game of Go with deep neural networks and tree search" (Google Deepmind team), Nature, January 2016

AlphaGoZero

No training on human expert knowledge – eventually becomes even better!



AlphaGoZero

Ke Jie stated that "After humanity spent thousands of years improving our tactics, computers tell us that humans are completely wrong... I would go as far as to say not a single human has touched the edge of the truth of Go."



Q-learning

An alternative to the policy gradient approach

Introduce a quality function Q that predicts the future reward for a given state s and a given action a. **Deterministic policy**: just select the action with the largest Q!



"value" of a state as color



"quality" of the action "going up" as color

Q-learning

Introduce a quality function Q that predicts the future reward for a given state s and a given action a. **Deterministic policy**: just select the action with the largest $\overline{Q}!$

$$Q(s_t, a_t) = E[R_t | s_t, a_t] \quad \text{(assuming future steps to follow the policy!)}$$

"Discounted"
future reward: $R_t = \sum_{t'=t}^T r_{t'} \gamma^{t'-t} \quad \text{depends on state and action at time t}$

Discount factor: $0 < \gamma \leq 1$

futu

ction at time t learning somewhat easier for smaller factor (short memory times)

Note: The 'value' of a state is $V(s) = \max_a Q(s, a)$ How do we obtain Q?

Q-learning: Update rule

Bellmann equation: (from optimal control theory) $Q(s_t, a_t) = E[r_t + \gamma \max_a Q(s_{t+1}, a) | s_t, a_t]$

In practice, we do not know the Q function yet, so we cannot directly use the Bellmann equation. However, the following update rule has the correct Q function as a fixed point:

$$Q^{\text{new}}(s_t, a_t) = Q^{\text{old}}(s_t, a_t) + \alpha (r_t + \gamma \max_a Q^{\text{old}}(s_{t+1}, a) - Q^{\text{old}}(s_t, a_t))$$

will be zero, once
we have converged
to the correct Q

If we use a neural network to calculate Q, it will be trained to yield the "new" value in each step.

Q(a=up,s)

Q(a=up,s)

Q(a=up,s)

Q-learning: Exploration

Initially, Q is arbitrary. It will be bad to follow this Q all the time. Therefore, introduce probability ϵ of random action ("exploration")!

Follow Q:"**exploitation**" Do something random (new):"**exploration**"

" ϵ -greedy"

Reduce this randomness later!

Example: Learning to play Atari Video Games

"Human-level control through deep reinforcement learning", Mnih et al., Nature, February 2015





last four 84x84 pixel images as input [=state]
motion as output [=action]

Example: Learning to play Atari Video Games

"Human-level control through deep reinforcement learning", Mnih et al., Nature, February 2015



Example: Learning to play Atari Video Games

"Human-level control through deep reinforcement learning", Mnih et al., Nature, February 2015



Apply RL to solve the challenge of finding, as fast as possible, a "treasure" in:

- a fixed given labyrinth
- an arbitrary labyrinth (in each run, the player finds itself in another labyrinth)

Use the labyrinth generator on Wikipedia "Maze Generation Algorithm"

Wikipedia: "Maze Generation Algorithm / Python Code Example"

